

**НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
ТЕХНОЛОГИЧЕСКИЙ УНИВЕРСИТЕТ
« М И С и С »
НОВОТРОИЦКИЙ ФИЛИАЛ**



Д.Д. Изаак

А.В.Швалёва

МАТЕМАТИЧЕСКАЯ СТАТИСТИКА

Лабораторный практикум

МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ РОССИЙСКОЙ ФЕДЕРАЦИИ
Федеральное государственное автономное образовательное учреждение
высшего профессионального образования
НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ ТЕХНОЛОГИЧЕСКИЙ УНИВЕРСИТЕТ
«МИСиС»

НОВОТРОИЦКИЙ ФИЛИАЛ

Кафедра математики и естествознания

Д.Д. Изаак
А.В. Швалёва

Математическая статистика

Лабораторный практикум

Новотроицк, 2012

УДК 519.25

ББК 22.17

ИЗ2

Научный редактор

Бонди И.Л., кандидат физико-математических наук

Рецензенты:

*Соколов А.А., кандидат физико-математических наук,
доцент кафедры общеобразовательных и профессиональных дисциплин
Орского филиала ФГАОУ ВПО
«Самарский государственный университет путей сообщения»*

*Попов А.С., кандидат педагогических наук,
доцент кафедры математического анализа, информатики,
теории и методики обучения информатики
Орского гуманитарно-технологического университета (филиала) ФГБОУ ВПО
«Оренбургский государственный университет»*

Изаак, Д. Д. Математическая статистика: лабораторный практикум /
Д. Д. Изаак, А. В. Швалёва – Магнитогорск: Издательский центр ФГБОУ ВПО
«МГТУ», 2012. – 51 с.

ISBN

Лабораторный практикум предназначен для студентов дневной и заочной форм обучения, изучающих курс «Теория вероятностей и математическая статистика». Разработан для студентов технических специальностей. Практикум предназначен для изучения возможностей программ MathCad, Excel, StatGraph и Stadia при обработке экспериментальных данных.

Рекомендовано Методическим советом НФ НИТУ «МИСиС»

ISBN	© Новотроицкий филиал ФГАОУ ВПО «Национальный исследовательский технологический университет "МИСиС", 2012
	© Изаак Д.Д., 2012
	© Швалёва А.В., 2012
	© Магнитогорский государственный технический университет им. Г.И. Носова, 2012

Содержание

Введение.....	4
Лабораторная работа 1. Непрерывные распределения.....	5
Лабораторная работа 2. Сравнение двух выборок.....	9
Лабораторная работа 3. Регрессионный анализ.....	27
Лабораторная работа 4. Исследование линейной корреляции.....	46
Библиографический список.....	51

Введение

Математическая статистика – это раздел математики, посвященный методам сбора, анализа и обработки статистических данных для научных и практических целей. Статистические данные представляют собой данные, полученные в результате обследования большого числа объектов и явлений. Обработка эмпирических данных, их систематизация, наглядное представление в форме графиков и таблиц, количественное описание посредством основных статистических показателей, формулировка выводов, имеющих прикладное значение для конкретной области человеческой деятельности – все это относится к методам математической статистики.

Однако, обрабатывая экспериментальные данные, приходится проводить очень трудоемкие вычисления. С появлением компьютеров такие вычисления стало проводить намного проще. В связи с этим будущим инженерам необходимо уметь проводить статистические расчеты не только аналитически, но и с помощью различных компьютерных программ. Данный лабораторный практикум посвящен изучению обработки экспериментальных данных с помощью программ MathCad, Excel, StatGraph и Stadia. В предлагаемом практикуме рассматриваются следующие разделы.

Раздел №1. Непрерывные случайные величины.

Раздел №2. Описательная статистика и сравнение двух выборок.

Раздел №3. Регрессионный анализ (однофакторный).

Раздел №4. Исследование линейной корреляции.

Каждая из предложенных четырех лабораторных работ содержит необходимые теоретические сведения, пояснения к работе с программами, а также разбор нулевого варианта. Лабораторные работы разработаны для десяти вариантов.

Данный лабораторный практикум предназначен для студентов технических направлений очной и заочной форм обучения.

Лабораторная работа №1: Непрерывные распределения

Используемое ПО: MathCad.

Цель работы: Научиться с помощью программы MathCad находить основные характеристики непрерывного распределения.

Задание

1. Построить график плотности вероятности.
2. Проверить, выполняется ли условие нормировки.
3. Найти функцию распределения и построить ее график.
4. Найти математическое ожидание, дисперсию и среднее квадратичное отклонение случайной величины.
5. Найти вероятность попадания случайной величины в указанный интервал с помощью плотности вероятности и с помощью функции распределения.
6. Найти медиану и квантиль, соответствующую указанной вероятности.

Таблица 1 – Содержание вариантов к лабораторной работе №1

№	Содержание варианта	№	Содержание варианта
1	$\varphi(x) = \begin{cases} 0, & x < 0 \\ \frac{1}{9}x^2, & 0 \leq x \leq 3 \\ 0, & x > 3 \end{cases}$ <p style="text-align: center;">(1,2) $p=0,1$</p>	2	$\varphi(x) = \begin{cases} 0, & x < 1 \\ \frac{4}{x^5}, & x \geq 1 \end{cases}$ <p style="text-align: center;">(2,3) $p=0,2$</p>
3	$\varphi(x) = \begin{cases} 0, & x < 1 \\ \frac{9}{x^{10}}, & x \geq 1 \end{cases}$ <p style="text-align: center;">(3,4) $p=0,3$</p>	4	$\varphi(x) = \begin{cases} 0, & x < 0 \\ e^{-x}, & x \geq 0 \end{cases}$ <p style="text-align: center;">(1,2) $p=0,4$</p>
5	$\varphi(x) = \begin{cases} 0, & x < 0 \\ 2 \cos 2x, & 0 \leq x \leq \frac{\pi}{4} \\ 0, & x > \frac{\pi}{4} \end{cases}$ <p style="text-align: center;">(4,5) $p=0,6$</p>	6	$\varphi(x) = \begin{cases} 0, & x < 0 \\ \frac{2}{9}(3x - x^2), & 0 \leq x \leq 3 \\ 0, & x > 3 \end{cases}$ <p style="text-align: center;">(1,2) $p=0,7$</p>

Окончание таблицы №1			
7	$\varphi(x) = 0,5e^{- x }$ $(1,5) \quad p=0,8$	8	$\varphi(x) = \begin{cases} 0, & x < 0 \\ 0,5 \sin x, & 0 \leq x \leq \pi \\ 0, & x > \pi \end{cases}$ $(1,2) \quad p=0,1$
9	$\varphi(x) = \begin{cases} 0, & x < -1 \\ \frac{2}{7}(2- x), & -1 \leq x \leq 2 \\ 0, & x > 2 \end{cases}$ $(1,2) \quad p=0,2$	10	$\varphi(x) = \begin{cases} 0, & x < 0 \\ 5e^{-5x}, & x \geq 0 \end{cases}$ $(3,6) \quad p=0,3$

Некоторые теоретические сведения

1-2. Функция плотности вероятности непрерывной случайной величины считается заданной корректно, если выполняются два условия:

а) $\varphi(x) \geq 0$;

б) $\int_{-\infty}^{\infty} \varphi(x) dx = 1$ (условие нормировки).

3. Функция распределения $F(x)$ случайной величины X определяется как

$$F(x) = \int_{-\infty}^x \varphi(t) dt .$$

4. Если случайная величина X имеет непрерывное распределение, то для любой функции $g(x)$

$$M(g(X)) = \int_{-\infty}^{\infty} g(x)\varphi(x)dx$$

(при условии, что интеграл сходится абсолютно).

Таким образом,

$$M(X) = \int_{-\infty}^{\infty} x\varphi(x)dx ,$$

$$D(X) = M(X^2) - M^2(X) = \int_{-\infty}^{\infty} x^2\varphi(x)dx - M^2(X) ,$$

$$\sigma(X) = \sqrt{D(X)} .$$

5. Вероятность попадания случайной величины в интервал есть интеграл от плотности вероятности по этому промежутку:

$$p(a < X < b) = \int_a^b \varphi(x) dx,$$

или приращение функции распределения на этом промежутке:

$$p(a < X < b) = F(b) - F(a).$$

6. Квантиль – это функция, обратная функции распределения. То есть, если $p = F(x)$, то квантиль $x_p = F^{-1}(p)$. В частном случае, когда $p = 0,5$, квантиль называют медианой. Таким образом, для нахождения квантили, соответствующей вероятности p , следует решить уравнение

$$F(x) - p = 0.$$

Пояснения к работе с программой

1. Для выполнения работы понадобятся панели инструментов «Графики», «Матанализ», «Арифметика», «Греческий алфавит», «Программирование», «Булево». Вывести их на экран можно через пункт меню «Вид / Панель инструментов».

2. Программу желательно составлять так, чтобы она обладала универсальностью. В том случае, если придется изменить данные: функцию плотности, интервал или вероятность, то пусть в программе это придется сделать всего один раз – в самом ее начале.

3. MathCad различает большие и маленькие буквы.

4. Задать функцию кусочно можно с помощью команд панели программирования «Add Line» и «If».

5. Построить декартов график можно с помощью команды панели графиков «Декартов график». При этом следует указать в нижнем поле ввода имя независимой переменной и в левом – функцию.

6. Если известно, что корень уравнения $f(x) = 0$ находится на отрезке $[a, b]$, то его можно найти командой **root**. Первый ее параметр – функция $f(x)$, второй – имя независимой переменной, в нашем случае x , третий и четвертый параметры – границы отрезка a и b . Так, например, один из корней уравнения $x^2 = 3$ можно найти так: **root(x²-3,x,1,2)**.

Программа MathCad

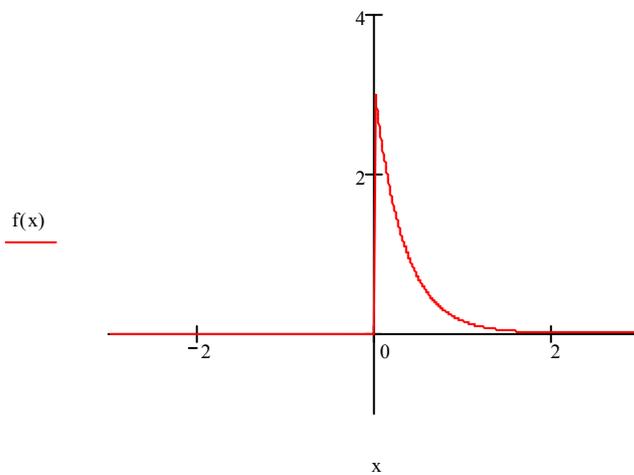
Разбор нулевого варианта

Пусть $\varphi(x) = \begin{cases} 3e^{-3x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$, интервал (1,2), вероятность $p=0,1$.

0. Условие задачи

$$f(x) := \begin{cases} 0 & \text{if } x < 0 \\ 3 \cdot e^{-3x} & \text{if } x \geq 0 \end{cases} \quad a := 1 \quad b := 2 \quad p := 0.1$$

1. График функции плотности

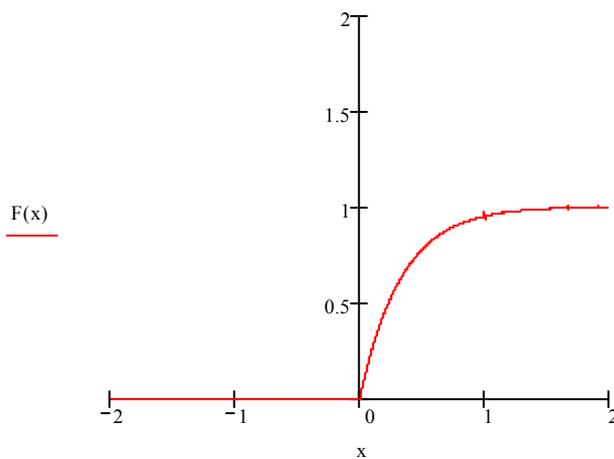


2. Условие нормировки

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

3. Функции распределения и ее график

$$F(x) := \int_{-\infty}^x f(t) dt$$



4. Математическое ожидание, дисперсия и стандартное отклонение

$$M := \int_{-\infty}^{\infty} x \cdot f(x) dx \quad M = 0.333$$

$$D := \int_{-\infty}^{\infty} x^2 \cdot f(x) dx - M^2 \quad D = 0.111$$

$$\sigma := \sqrt{D} \quad \sigma = 0.333$$

5. Вероятность попадания в интервал

$$p1 := \int_a^b f(x) dx \quad p1 = 0.047$$

$$p2 := F(b) - F(a) \quad p2 = 0.047$$

6. Медиана и квантиль

$$\text{root}(F(x) - 0.5, x, 0, 2) = 0.231$$

$$\text{root}(F(x) - p, x, 0, 2) = 0.035$$

Лабораторная работа №2: Сравнение двух выборок

Используемое ПО: MathCad, Excel, StatGraph, Stadia.

Цель работы: Научиться с помощью вышеуказанных программ находить основные характеристики случайных выборок и сравнивать их.

Задание

1. Найти средние арифметические и эмпирические стандарты для каждой из выборок.
2. Построить доверительные интервалы для математических ожиданий и стандартных отклонений.
3. Проверить гипотезу о равенстве дисперсий.
4. Если гипотеза о равенстве дисперсий принята, найти сводную оценку стандартного отклонения.
5. Проверить гипотезу о равенстве математических ожиданий.
6. Если гипотеза о равенстве математических ожиданий принята, найти сводную оценку математического ожидания и объединенную оценку стандартного отклонения.
7. Объединить две выборки в одну и проверить гипотезу о том, что экспериментальные данные имеют нормальный закон распределения. Рассматривать интервалы равной длины. Число интервалов равно L .
8. Построить гистограмму.

Указания:

- 1) Во всех пунктах брать уровень значимости $\alpha = 0,05$.
- 2) Выполнить пункты 1-7 в программе MathCad. В пункте 7 искать только теоретическую квантиль.
- 3) Выполнить пункты 1,2,3,5,7 в программе Excel. Доверительные интервалы для математических ожиданий искать с помощью встроенной команды и без нее. В пункте 7 искать только теоретическую квантиль.
- 4) Выполнить пункты 1,2,3,5,7,8 в программе StatGraph.
- 5) Выполнить пункты 1,2,3,5,7,8 в программе Stadia.

Содержание вариантов к лабораторной работе № 2

Вариант 1. $L=7$.

1 серия измерений. $n_1 = 28$.

7,2	7,1	3,7	5,3	6,4	5,2	9,7	8,8	6,3
5,9	6,9	4,5	9,0	5,5	6,1	7,5	8,9	3,7
6,3	6,1	6,3	7,2	6,2	3,5	9,0	6,4	7,5
9,8								

2 серия измерений. $n_2 = 31$.

6,7	5,1	7,4	5,9	9,8	6,6	8,8	9,3	7,9
5,6	7,2	6,2	6,8	5,4	6,8	8,2	9,3	8,0
6,0	6,0	7,6	7,5	8,9	4,9	5,8	8,5	8,9
8,7	6,4	6,6	5,7					

Вариант 2. $L=7$.

1 серия измерений $n_1 = 38$.

23,2	20,1	18,8	24,1	21,6	22,8	22,1	25,2	24,8
20,6	25,9	27,0	21,9	23,5	21,7	21,1	21,3	18,3
21,0	23,8	17,4	17,3	17,9	20,6	18,4	24,2	20,7
22,0	18,3	22,6	20,2	21,5	16,5	21,3	21,5	17,9
26,2	29,1							

2 серия измерений. $n_2 = 21$.

24,2	30,4	21,7	21,2	20,7	25,4	17,8	19,0	21,5
18,7	22,1	25,6	19,6	19,9	21,9	26,4	21,9	24,1
21,1	23,5	23,9						

Вариант 3. $L=8$.

1 серия измерений. $n_1 = 36$.

4,4	4,3	2,9	3,8	5,1	5,6	3,7	4,7	5,6
4,4	4,6	4,0	4,7	5,0	5,4	4,3	6,2	6,4
4,5	5,1	4,4	5,4	5,8	4,8	6,2	5,0	5,6
4,8	4,7	4,0	5,3	5,4	2,5	5,4	5,4	6,9

2 серия измерений. $n_2 = 22$.

5,0	4,3	3,9	6,0	3,4	4,1	4,7	3,4	4,2
5,1	3,5	3,1	4,3	3,7	3,7	6,2	4,8	3,5
4,3	6,2	2,7	7,1					

Вариант 4. L=7.

1 серия измерений $n_1 = 34$.

0,37	0,67	0,64	0,94	0,82	0,60	0,70	0,84	0,81
0,67	0,77	0,52	0,70	0,54	0,76	0,47	0,86	0,62
0,54	0,88	0,85	0,84	0,62	0,65	0,70	0,97	0,55
0,72	0,74	0,93	0,65	0,61	0,66	0,78		

2 серия измерений $n_2 = 23$.

0,84	0,85	0,84	0,64	0,97	0,78	0,81	0,73	0,93
0,60	0,79	0,69	0,85	0,73	0,72	0,84	0,54	0,95
0,42	0,90	0,90	0,87	0,78				

Вариант 5. L=8.

1 серия измерений. $n_1 = 32$.

23,5	26,5	22,8	26,0	23,7	25,8	22,7	18,9	24,6
26,1	26,2	20,9	23,0	31,9	23,2	21,3	21,8	24,1
26,7	19,3	27,1	26,6	29,5	20,5	21,8	26,8	23,9
20,3	21,5	24,8	25,8	20,5				

2 серия измерений. $n_2 = 24$.

24,6	20,4	20,1	25,7	24,0	24,2	19,5	18,3	21,2
18,8	15,4	22,1	24,2	20,4	19,2	15,2	18,0	21,5
22,2	21,4	17,0	17,8	25,4	21,5			

Вариант 6. L=6.

1 серия измерений. $n_1 = 30$.

4,9	4,8	3,6	4,8	6,5	8,6	5,5	7,7	6,1
7,0	7,5	7,9	5,7	6,2	6,6	5,0	7,2	5,9
7,2	4,0	6,3	6,1	5,0	6,7	3,6	6,4	3,9
3,5	5,4	6,5						

2 серия измерений. $n_2 = 25$.

4,5	5,8	5,0	7,1	5,7	6,2	6,2	6,6	5,6
4,9	5,5	5,7	5,3	6,7	5,5	7,1	6,2	3,0
6,3	7,9	5,9	9,5	6,7	7,3	5,5		

Вариант 7. L=7.

1 серия измерений. $n_1 = 28$.

0,22 0,75 0,77 0,59 0,85 0,99 0,79 1,10 1,04
0,51 0,83 0,72 0,65 1,03 0,92 0,53 0,63 0,85
0,91 0,72 0,62 0,58 0,81 0,91 0,81 1,02 1,15
0,35

2 серия измерений. $n_2 = 26$.

0,65 0,96 0,85 0,98 0,69 1,01 0,79 0,62 0,71
0,84 1,13 0,81 1,02 0,65 0,54 0,78 0,69 0,65
0,61 0,61 0,72 1,01 0,54 0,58 0,70 0,82

Вариант 8. L=6.

1 серия измерений. $n_1 = 26$.

29,8 29,5 30,4 30,4 28,5 35,6 29,3 28,0 26,4
24,2 32,3 26,2 22,9 25,9 27,5 20,2 28,4 22,7
21,3 23,3 23,2 29,7 24,0 26,5 28,5 24,7

2 серия измерений. $n_2 = 27$.

27,9 32,9 29,1 31,1 28,9 34,1 35,7 23,7 33,9
25,2 25,3 29,3 29,7 29,9 36,7 24,7 32,5 30,4
26,4 31,1 28,8 34,7 32,9 36,1 30,5 39,7 30,8

Вариант 9. L=7.

1 серия измерений. $n_1 = 24$.

6,7 6,6 5,7 5,8 6,5 6,5 7,5 5,6 6,3
5,9 3,9 5,8 8,0 7,3 7,9 5,5 9,4 6,3
6,2 7,8 7,1 9,4 7,6 7,3

2 серия измерений. $n_2 = 28$.

4,9 6,0 3,7 6,6 5,5 4,7 8,0 6,7 4,7
4,4 5,3 5,7 5,1 5,7 6,5 6,3 6,9 5,1
6,2 4,2 2,0 5,7 5,6 7,3 5,1 6,6 6,1
5,6

Вариант 10. $L=6$.

1 серия измерений. $n_1 = 22$.

0,82 0,61 1,10 0,51 0,44 0,65 0,88 0,58 0,68
0,89 1,07 1,15 0,96 1,01 0,49 0,99 0,90 1,10
0,74 0,88 1,09 1,22

2 серия измерений. $n_2 = 29$.

1,12 0,90 1,03 1,35 1,38 0,77 1,05 0,90 0,60
0,74 1,04 1,00 1,32 0,52 1,13 0,68 0,90 1,04
0,66 0,95 0,66 0,99 0,95 1,19 0,90 1,26 1,12
0,99 1,14

Некоторые теоретические сведения

1. Несмещенными оценками математического ожидания и стандартного отклонения являются среднее арифметическое и эмпирический стандарт:

$$\bar{y} = \sum_{i=1}^n \frac{y_i}{n}, \quad s = \sqrt{\sum_{i=1}^n \frac{(y_i - \bar{y})^2}{n-1}}.$$

2. Доверительные интервалы для математического ожидания и стандартного отклонения соответственно имеют вид:

$$\bar{y} - |t|_p \frac{s}{\sqrt{n}} < \beta < \bar{y} + |t|_p \frac{s}{\sqrt{n}}, \quad (2.1)$$

$$s \cdot \sqrt{\frac{k}{a_2}} < \sigma < s \cdot \sqrt{\frac{k}{a_1}}.$$

Здесь $p = 1 - \alpha$, $k = n - 1$, a_1 и a_2 – квантили распределения Пирсона для вероятностей $\frac{1-p}{2}$ и $\frac{1+p}{2}$ соответственно. Видно, что доверительный интервал для математического ожидания симметричен относительно среднего арифметического, поэтому иногда считают, что интервал найден, если указаны среднее арифметическое и радиус интервала $R = |t|_p \frac{s}{\sqrt{n}}$.

3. Пусть $F_{\mathcal{D}_1}$ – отношение большей эмпирической дисперсии к меньшей, $F_{\mathcal{D}_2}$ – меньшей к большей. Пусть F_{m_1} и F_{m_2} – квантили распределения Фишера $F(p_1, k_1, k_2)$ и $F(p_2, k_2, k_1)$. Здесь $p_1 = 1 - \frac{\alpha}{2}$, $p_2 = \frac{\alpha}{2}$, k_1 – число степеней свободы большей эмпирической дисперсии k_2 – меньшей. Критерий Фишера

требует отвергнуть гипотезу о равенстве дисперсий, если $F_{\mathcal{E}1} > F_{m1}$ ($F_{\mathcal{E}2} < F_{m2}$). В противном случае гипотеза принимается.

4. Если гипотеза о равенстве дисперсий не противоречит экспериментальным данным, можно найти сводную оценку дисперсии

$$s_{св}^2 = \frac{s_1^2 k_1 + s_2^2 k_2}{k_1 + k_2},$$

где s_1^2 и s_2^2 – эмпирические дисперсии с числом степеней свободы k_1 и k_2 соответственно.

5. Пусть

$$T_{\mathcal{E}} = \frac{|\bar{y}_1 - \bar{y}_2|}{s_{св} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}},$$

а $|t|_p$ – квантиль распределения модуля отношения Стьюдента для вероятности $p = 1 - \alpha$ и числа степеней свободы $k = k_1 + k_2$. Критерий Стьюдента требует отвергнуть гипотезу о равенстве математических ожиданий, если $T_{\mathcal{E}} > |t|_p$. В противном случае гипотеза принимается.

6. Если гипотеза о равенстве математических ожиданий не противоречит экспериментальным данным, можно найти сводную оценку математического ожидания и объединенную оценку дисперсии:

$$\bar{y}_{св} = \frac{\bar{y}_1 n_1 + \bar{y}_2 n_2}{n_1 + n_2}, \quad s_{об}^2 = \frac{1}{n-1} \left(\sum_{i=1}^{n_1} y_i^2 + \sum_{i=1}^{n_2} y_i'^2 - n \cdot \bar{y}_{св}^2 \right).$$

Здесь $n = n_1 + n_2$, $\{y_i\}$ и $\{y_i'\}$ – данные выборки. Заметим, что $\bar{y}_{св}$ и $s_{об}$ есть обычное среднее арифметическое и обычный эмпирический стандарт для объединенной выборки.

7. Для проверки гипотезы о том, что экспериментальные данные распределены по нормальному закону, можно применить критерий согласия Пирсона. Для этого весь диапазон изменения случайной величины разбивают на L участков точками $x_1 = -\infty$, x_2 , ..., $x_{L+1} = +\infty$, подсчитывают частоты попадания в эти участки $\frac{N_1}{n}$, $\frac{N_2}{n}$, ..., $\frac{N_L}{n}$, а также теоретические вероятности попадания в эти же участки:

$$p_i = \frac{1}{\sqrt{2\pi}} \int_{x_i}^{x_{i+1}} e^{-\frac{U^2}{2}} dU.$$

Затем вычисляют значение величины

$$\chi^2 = \sum_{i=1}^L \frac{(N_i - np_i)^2}{np_i}$$

и сравнивают его с квантилью распределения Пирсона $\chi_p^2(k)$ для вероятности $p = 1 - \alpha$ и числа степеней свободы $k = L - 3$. Критерий Пирсона требует отвергнуть гипотезу о том, что экспериментальные данные распределены по нормальному закону, если $\chi^2 > \chi_p^2(k)$, в противном случае гипотеза принимается. На практике точки дробления x_i выбирают так, чтобы либо все интервалы имели одинаковую длину, либо все вероятности p_i равнялись бы между собой.

8. Если по оси абсцисс отметить точки x_i из пункта 7, только вместо x_1 и x_{L+1} взять наименьшее и наибольшее числа в выборке, и на каждом отрезке построить прямоугольник высотой $\frac{N_i}{n \cdot h_i}$, где h_i – длина i -го участка, то графическое изображение полученных данных даст гистограмму.

Пояснения к работе с программой MathCad

1. Для выполнения работы понадобятся панели инструментов «Арифметика», «Греческий алфавит», «Программирование», «Матрицы». Вывести их на экран можно через пункт меню «Вид / Панель инструментов».

2. Программу желательно составлять так, чтобы она обладала универсальностью. В том случае, если придется изменить выборки, то пусть в программе при этом ничего не придется менять.

3. Перед тем, как приступить к работе с MathCad'ом, следует подготовить два текстовых файла с выборками, например, с помощью программы Блокнот. Набирать данные следует либо в одну строку, разделяя их пробелами, либо в один столбец. В качестве десятичного разделителя следует использовать точку. Сохраните эти текстовые файлы в той же папке, где будет храниться файл, созданный в MathCad'е.

4. Укажем некоторые команды, которые нам понадобятся для выполнения работы.

READPRN("FileName.txt")

Считать данные из текстового файла.

length(x)

Длина массива x .

mean(x)

Среднее арифметическое массива x .

stdev(x)

Смещенная оценка стандартного отклонения $\tilde{s} = \sqrt{\sum_{i=1}^n \frac{(y_i - \bar{y})^2}{n}}$.

qchisq(p,k)

Квантиль распределения Пирсона для вероятности p и числа степеней свободы k .

qt(p,k)

Обычная квантиль (не модуля) отношения Стьюдента для вероятности p и числа степеней свободы k . Заметим, что из обычной квантили $qt(p,k)$ можно получить нужную нам квантиль модуля так: $-qt\left(\frac{1-p}{2}, k\right)$.

qF(p,k₁,k₂)

Квантиль распределения Фишера для вероятности p и числа степеней свободы k_1 и k_2 .

stack(x,y)

Объединение двух массивов-столбцов в один.

augment(x,y)

Объединение двух массивов-строк в один.

x:=x^T

Транспонирование массива x . Некоторые команды программы MathCad применимы к массивам-строкам, некоторые к массивам-столбцам. С помощью команды транспонирования (панель «*Матрицы*») можно преобразовывать рассматриваемый массив.

Пояснения к работе с программой Excel

1. Список всех доступных команд можно получить с помощью кнопки быстрого доступа «*Вставка функции*» или через меню «*Вставка / Функция...*».
2. Укажем некоторые команды, которые нам понадобятся для выполнения работы.

СРЗНАЧ(A_i:B_j)

Среднее арифметическое диапазона ячеек $A_i:B_j$.

СТАНДОТКЛОН($A_i:B_j$)

Эмпирический стандарт диапазона ячеек $A_i:B_j$.

КОРЕНЬ(A_i)

Квадратный корень.

ХИ2ОБР(α, k)

Квантиль распределения Пирсона для уровня значимости α и числа степеней свободы k .

СТЬЮДРАСПОБР(α, k)

Квантиль распределения модуля отношения Стьюдента для уровня значимости α и числа степеней свободы k .

ФРАСПОБР(α, k_1, k_2)

Квантиль распределения Фишера для уровня значимости α и числа степеней свободы k_1 и k_2 .

ДОВЕРИТ(α, s, n)

Радиус доверительного интервала для математического ожидания. Здесь α – уровень значимости, s – эмпирический стандарт, n – число элементов в выборке. Заметим, что эта команда находит радиус, заменяя в формуле (2.1) квантиль модуля отношения Стьюдента квантилью модуля стандартного нормального распределения. Учитывая то, что при $n \rightarrow \infty$ распределение Стьюдента стремится к стандартному нормальному, при достаточно больших n эти квантили практически совпадают.

Пояснения к работе с программой StatGraph

1. Следует подготовить три столбца с экспериментальными данными: два с соответствующими выборками и третий – с объединенной выборкой.

2. В качестве десятичного разделителя следует использовать запятую.

3. Для выполнения пунктов 1,2,3,5 используйте пункт меню «*Compare / Two Samples / Two-sample Comparison...*». В окошке «*Two-sample Comparison*» используйте кнопку быстрого доступа «*Tabular Options*» для получения доступа к разделам «*Comparison Of Means*» и «*Comparison Of Standard Deviations*».

4. Для выполнения пунктов 7,8 используйте пункт меню «*Describe / Distribution Fitting / Uncensored Data...*». В окошке «*Uncensored Data*» используйте кнопки быстрого доступа «*Tabular Options*» и «*Graphical Options*» для получения доступа к разделам «*Test For Normality*» и «*Frequency Histogram*» соответственно.

Пояснения к работе с программой Stadia

1. Следует подготовить три столбца с экспериментальными данными: два с соответствующими выборками и третий – с объединенной выборкой.
2. В качестве десятичного разделителя следует использовать точку.
3. Для выполнения всех пунктов воспользуемся пунктом меню «Статист-F9».
4. Для выполнения пунктов 1,2 воспользуемся кнопкой «Описательная статистика».
5. Для выполнения пунктов 3,5 воспользуемся кнопкой «Стьюдента и Фишера».
6. Для выполнения пунктов 7,8 воспользуемся кнопкой «Гистограмма/Нормальность».

Разбор нулевого варианта

Вариант 0. $L=8$.

1 серия измерений. $n_1 = 26$.

0,44 1,29 1,25 1,06 1,24 0,87 0,96 1,25 0,88
1,09 0,90 0,75 1,09 0,73 1,04 1,07 1,20 1,15
0,94 0,83 1,19 0,88 0,77 0,84 0,79 0,80

2 серия измерений. $n_2 = 32$.

0,89 0,97 0,33 0,93 0,84 1,43 1,34 0,88 0,75
0,96 1,09 0,83 0,95 0,33 0,56 1,20 1,12 0,92
0,73 1,30 0,70 1,27 0,82 0,86 1,30 1,00 1,12
0,69 1,03 0,58 1,26 1,16

Программа MathCad

Разбор нулевого варианта

0. Подготовка

X := READPRN("00Vib1.txt") Y := READPRN("00Vib2.txt")

n1 := length(X) n2 := length(Y)

k1 := n1 - 1 k2 := n2 - 1

a := 0.05 p := 1 - a

1. Средние и стандарты

x1 := mean(X) x1 = 0.973

y1 := mean(Y) y1 = 0.942

s1 := stdev(X) · $\sqrt{\frac{n1}{k1}}$ s1 = 0.207

s2 := stdev(Y) · $\sqrt{\frac{n2}{k2}}$ s2 = 0.276

	0
0	0.44
1	1.29
2	1.25
3	1.06
4	1.24
5	0.87
6	0.96
7	1.25
8	0.88
9	1.09
10	0.9
11	0.75
12	1.09
13	0.73
14	1.04
15	1.07

X =

	0
0	0.89
1	0.97
2	0.33
3	0.93
4	0.84
5	1.43
6	1.34
7	0.88
8	0.75
9	0.96
10	1.09
11	0.83
12	0.95
13	0.33
14	0.56
15	1.2

Y =

2. Доверительные интервалы

Для сигма 1

$p1 := \frac{1-p}{2}$ $p2 := \frac{1+p}{2}$ a1 := qchisq(p1, k1) a2 := qchisq(p2, k1)

sig11 := s1 · $\sqrt{\frac{k1}{a1}}$ sig12 := s1 · $\sqrt{\frac{k1}{a2}}$ sig11 = 0.163 sig12 = 0.286

Для сигма 2

a1 := qchisq(p1, k2) a2 := qchisq(p2, k2)

sig21 := s2 · $\sqrt{\frac{k2}{a1}}$ sig22 := s2 · $\sqrt{\frac{k2}{a2}}$ sig21 = 0.221 sig22 = 0.366

Для бета 1

$tp := -qt\left(\frac{1-p}{2}, k1\right)$ r1 := tp · $\frac{s1}{\sqrt{n1}}$ beta11 := x1 - r1 beta11 = 0.889

beta12 := x1 + r1 beta12 = 1.057

Для бета 2

$tp := -qt\left(\frac{1-p}{2}, k2\right)$ r2 := tp · $\frac{s2}{\sqrt{n2}}$ beta21 := y1 - r2 beta21 = 0.843

beta22 := y1 + r2 beta22 = 1.041

3. Гипотеза о равенстве дисперсий

$$F1 := \begin{cases} \frac{s1^2}{s2^2} & \text{if } (s1 > s2) \\ \frac{s2^2}{s1^2} & \text{if } (s2 \geq s1) \end{cases} \quad K1 := \begin{cases} k1 & \text{if } (s1 > s2) \\ k2 & \text{if } (s2 \geq s1) \end{cases} \quad K2 := \begin{cases} k1 & \text{if } (s2 > s1) \\ k2 & \text{if } (s1 \geq s2) \end{cases}$$

$$F1 = 1.768 \quad F2 := qF\left(1 - \frac{\alpha}{2}, K1, K2\right) \quad F2 = 2.174$$

$$Flag := \begin{cases} \text{"No"} & \text{if } (F1 > F2) \\ \text{"Yes"} & \text{if } (F2 \geq F1) \end{cases} \quad Flag = \text{"Yes"}$$

4. Сводная оценка стандартного отклонения

$$Ssv := \frac{s1^2 \cdot k1 + s2^2 \cdot k2}{k1 + k2} \quad Ssv := \sqrt{Ssv} \quad Ssv = 0.247$$

5. Гипотеза о равенстве математических ожиданий

$$T1 := \frac{|x1 - y1|}{Ssv \cdot \sqrt{\frac{1}{n1} + \frac{1}{n2}}} \quad T1 = 0.478 \quad T2 := -qt\left(\frac{1 - p}{2}, k1 + k2\right) \quad T2 = 2.003$$

$$Flag := \begin{cases} \text{"No"} & \text{if } (T1 > T2) \\ \text{"Yes"} & \text{if } (T2 \geq T1) \end{cases} \quad Flag = \text{"Yes"}$$

6. Сводная оценка математического ожидания, объединенный стандарт

$$Ysv := \frac{x1 \cdot n1 + y1 \cdot n2}{n1 + n2} \quad Ysv = 0.956 \quad z := \text{stack}(X, Y)$$

$$Sob := \text{stdev}(z) \cdot \sqrt{\frac{\text{length}(z)}{\text{length}(z) - 1}} \quad Sob = 0.246$$

7. Квантиль распределения Пирсона для критерия согласия

$$\chi := qchisq(p, 5) \quad \chi = 11.07$$

	0
0	0.44
1	1.29
2	1.25
3	1.06
4	1.24
5	0.87
6	0.96
z = 7	1.25
8	0.88
9	1.09
10	0.9
11	0.75
12	1.09
13	0.73
14	1.04
15	1.07

Программа Excel

Разбор нулевого варианта

0,44	0,89	26	x1-сред.	y1-сред.	S1	S2
1,29	0,97	32	0,973077	0,941875	0,207283	0,275616
1,25	0,33	8	Сигма1-1	Сигма1-2	Сигма2-1	Сигма2-2
1,06	0,93		0,162563	0,286135	0,220962	0,366426
1,24	0,84		Квантили ст. норм. распр.			
0,87	1,43		Бета1-1	Бета1-2	Бета2-1	Бета2-2
0,96	1,34		0,893402	1,052752	0,846381	1,037369
1,25	0,88		Квантили распр. Стьюд.			
0,88	0,75		Бета1-1	Бета1-2	Бета2-1	Бета2-2
1,09	0,96		0,889354	1,0568	0,842505	1,041245
0,9	1,09		F-эсп.	F-квант.	F-эсп.	F-квант.
0,75	0,83		0,565611	0,460064	1,767999	2,173607
1,09	0,95		T-эсп.	T-квант.		
0,73	0,33		0,47757	2,003239		
1,04	0,56		X2-квант.			
1,07	1,2		11,07048			
1,2	1,12					
1,15	0,92					
0,94	0,73					
0,83	1,3					
1,19	0,7					
0,88	1,27					
0,77	0,82					
0,84	0,86					
0,79	1,3					
0,8	1					
	1,12					
	0,69					
	1,03					
	0,58					
	1,26					
	1,16					

Программа StatGraph

Разбор нулевого варианта

Comparison of Means

95,0% confidence interval for mean of Col_1: 0,973077 +/- 0,0837235

95,0% confidence interval for mean of Col_2: 0,941875 +/- 0,0993703

95,0% confidence intervals for the difference between the means:

assuming equal variances: 0,0312019 +/- 0,130882

not assuming equal variances: 0,0312019 +/- 0,127129

t tests to compare means

Null hypothesis: mean1 = mean2

(1) Alt. hypothesis: mean1 NE mean2

assuming equal variances: t = 0,47757 P-value = 0,634815

not assuming equal variances: t = 0,491724 P-value = 0,624846

(2) Alt. hypothesis: mean1 > mean2

assuming equal variances: t = 0,47757 P-value = 0,317408

not assuming equal variances: t = 0,491724 P-value = 0,312423

(3) Alt. hypothesis: mean1 < mean2

assuming equal variances: t = 0,47757 P-value = 0,682592

not assuming equal variances: t = 0,491724 P-value = 0,687577

The StatAdvisor

This option runs a t-test to compare the means of the two samples. It also constructs confidence intervals for each mean and for the difference between the means. Of particular interest is the confidence interval for the difference between the means, which extends from -0,0996797 to 0,162084. Since the interval contains the value 0.0, there is not a statistically significant difference between the means of the two samples at the 95,0% confidence level. The t-tests can also be used to arrive at the same conclusion. P-values below 0,05 indicate significant differences between the two means. NOTE: the interval used above assumes that the variances of the two samples are equal. This was determined by running an F-test to compare the standard deviations of the two samples. You can see the results of that test by selecting Comparison of Standard Deviations from the Tabular Options menu.

Comparison of Standard Deviations

	Col_1	Col_2
Standard deviation	0,207283	0,275616
Variance	0,0429662	0,0759641
Df	25	31
Ratio of Variances = 0,565611		

95,0% Confidence Intervals

Standard deviation of Col_1: [0,162563;0,286135]
Standard deviation of Col_2: [0,220962;0,366426]
Ratio of Variances: [0,26847;1,22942]

F-tests to Compare Standard Deviations

Null hypothesis: $\sigma_1 = \sigma_2$

- (1) Alt. hypothesis: $\sigma_1 \neq \sigma_2$
F = 0,565611 P-value = 0,147687
- (2) Alt. hypothesis: $\sigma_1 > \sigma_2$
F = 0,565611 P-value = 0,926156
- (3) Alt. hypothesis: $\sigma_1 < \sigma_2$
F = 0,565611 P-value = 0,0738437

The StatAdvisor

This option runs an F-test to compare the variances of the two samples. It also constructs confidence intervals for each standard deviation and for the ratio of the variances. Of particular interest is the confidence interval for the ratio of the variances, which extends from 0,26847 to 1,22942. Since the interval contains the value 1.0, there is not a statistically significant difference between the standard deviations of the two samples at the 95,0% confidence level. The F-tests can also be used to arrive at the same conclusion. P-values below 0,05 indicate significant differences between the two standard deviations. IMPORTANT NOTE: the F-tests and confidence intervals shown here depend on the samples having come from normal distributions. To test this assumption, select Summary Statistics from the list of Tabular Options and check the standardized skewness and standardized kurtosis values.

Tests for Normality for Col_4

Computed Chi-Square goodness-of-fit statistic = 10,1379
P-Value = 0,859327

Shapiro-Wilks W statistic = 0,964943
P-Value = 0,193093

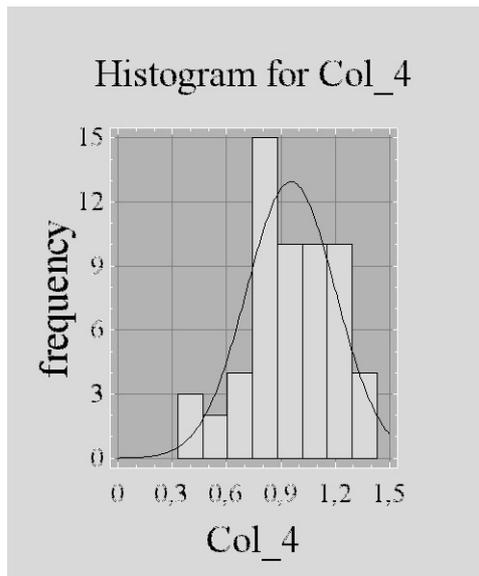
Z score for skewness = 0,992315
P-Value = 0,321042

Z score for kurtosis = 0,389011
P-Value = 0,697264

The StatAdvisor

This pane shows the results of several tests run to determine whether Col_4 can be adequately modeled by a normal distribution. The chi-square test divides the range of Col_4 into 19 equally probable classes and compares the number of observations in each class to the number expected. The Shapiro-Wilks test is based upon comparing the quantiles of the fitted normal distribution to the quantiles of the data. The standardized skewness test looks for lack of symmetry in the data. The standardized kurtosis test looks for distributional shape which is either flatter or more peaked than the normal distribution.

The lowest P-value amongst the tests performed equals 0,193093. Because the P-value for this test is greater than or equal to 0.10, we can not reject the idea that Col_4 comes from a normal distribution with 90% or higher confidence.



Программа Stadia

Разбор нулевого варианта

ОПИСАТЕЛЬНАЯ СТАТИСТИКА. Файл: var0.std

Переменная	Размер	<---Диапазон--->		Среднее---Ошибка		Дисперс	Ст.откл	Сумма
vib1	26	0,44	1,29	0,9731	0,04065	0,04297	0,2073	25,3
vib2	32	0,33	1,43	0,9419	0,04872	0,07596	0,2756	30,14

Переменная	Медиана	<--Квартили-->		ДовИнтСр.	<-ДовИнтДисп->		Ош.СтОткл
vib1	0,95	0,8225	1,16	0,08268	0,02643	0,0819	0,07776
vib2	0,94	0,7675	1,15	0,09813	0,04883	0,1343	0,09802

Переменная	Асимметр.	Значим	Эксцесс	Значим
vib1	-0,363	0,1986	2,812	0,4812
vib2	-0,3621	0,1795	2,747	0,4587

КРИТЕРИЙ ФИШЕРА И СТЬЮДЕНТА. Файл: var0.std

Переменные: vib1, vib2

Статистика Фишера=0,5656, Значимость=0,0737, степ.своб = 31,25

Гипотеза 0: <Нет различий между выборочными дисперсиями>

Статистика Стьюдента=0,492, Значимость=0,6302, степ.своб = 56

Гипотеза 0: <Нет различий между выборочными средними>

ГИСТОГРАММА И ТЕСТ НОРМАЛЬНОСТИ. Файл: var0.std

Х-лев.	Х-станд	Частота	%	Накопл.	%
0,33	-2,547	3	5,172	3	5,172
0,4675	-1,987	2	3,448	5	8,621
0,605	-1,428	4	6,897	9	15,52
0,7425	-0,8681	15	25,86	24	41,38
0,88	-0,3087	10	17,24	34	58,62
1,017	0,2508	10	17,24	44	75,86
1,155	0,8103	10	17,24	54	93,1
1,293	1,37	4	6,897	58	100
1,43	1,929				

Колмогоров=0,05836, Значимость=1,582, степ.своб = 58

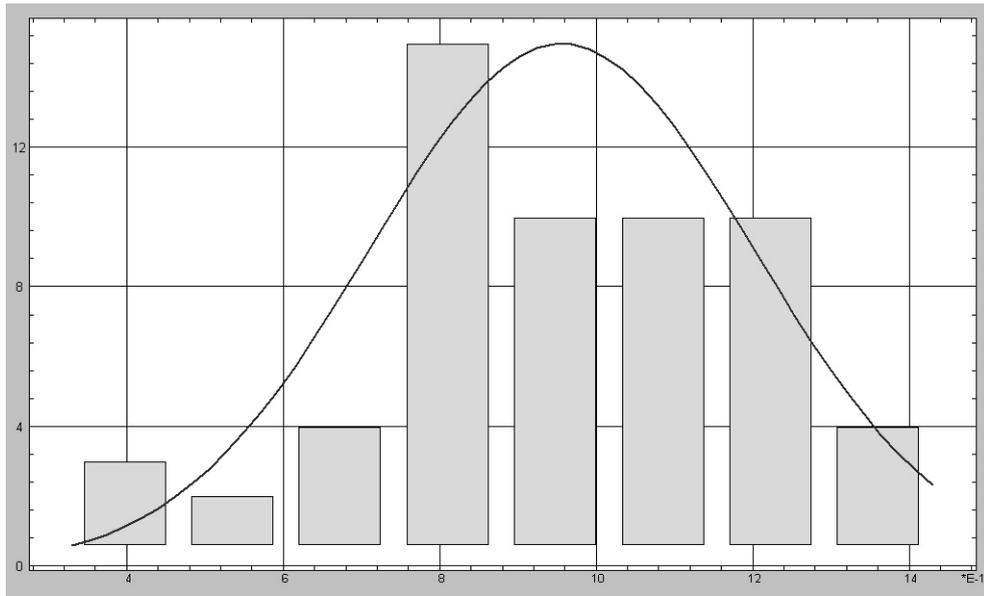
Гипотеза 0: <Распределение не отличается от нормального>

Омега-квадрат=0,04288, Значимость=0,6628, степ.своб = 58

Гипотеза 0: < Распределение не отличается от нормального>

Хи-квадрат=8,909, Значимость=0,1127, степ.своб = 5

Гипотеза 0: < Распределение не отличается от нормального>



Лабораторная работа №3: Регрессионный анализ

Используемое ПО: MathCad, Excel, StatGraph, Stadia.

Цель работы: Научиться с помощью вышеуказанных программ строить регрессионные модели, проверять их адекватность и значимость коэффициентов.

Задание

1. Найти линейное, квадратичное и кубическое регрессионные уравнения, уравнение вида $y = e^{ax+b}$. Изобразить на одном графике найденную экспоненциальную функцию и результаты эксперимента. Построить график остатков для экспоненциальной модели. Данные считать из файла. В качестве значений y брать средние арифметические. Все веса считать равными 1 (Программа MathCad).

2. Найти линейное уравнение регрессии. В качестве значений y брать средние арифметические. Все веса считать равными 1. Проверить адекватность построенной модели и значимость коэффициентов (Программа Excel).

3. Найти линейное и квадратичное уравнения регрессии с учетом весов. Для каждой регрессионной модели построить график найденной функции и результатов эксперимента. Проверить адекватность каждой построенной модели и значимость всех коэффициентов (Программа StatGraph).

4. Найти линейное, квадратичное и кубическое регрессионные уравнения, уравнение вида $y = e^{ax+b}$. Изобразить на одном графике найденную квадратичную функцию и результаты эксперимента. В качестве значений y брать средние арифметические. Все веса считать равными 1. Проверить адекватность каждой построенной модели и значимость всех коэффициентов (Программа Stadia).

5. Напишите программу в MathCad для построения полноценных линейной и квадратичной регрессионных моделей с учетом весов, проверки адекватности построенных моделей и значимости всех коэффициентов.

Таблица 2 – Содержание вариантов к лабораторной работе № 3

Вариант 1					
<i>X</i>	<i>Y</i>				
0,00	4,96	4,94	4,96	4,94	
0,10	5,42	5,46	5,42	5,39	
0,15	5,63	5,64	5,63	5,64	
0,25	6,00	6,03	6,02		
0,30	6,19	6,25	6,20		
0,55	7,04	7,04	7,04		
0,70	7,46	7,41	7,43		
0,75	7,54	7,54	7,53	7,60	
0,80	7,66	7,65	7,64	7,64	
0,90	7,87	7,79	7,81	7,86	
Вариант 2					
<i>X</i>	<i>Y</i>				
5	53	52	53		
10	63	63	62	64	
25	91	91	91	90	
30	96	94	97	97	
40	108	108	107		
55	117	118	116		
70	117	116	117	117	
75	114	114	114	115	
90	103	102	101	106	
100	89	89	91		
Вариант 3					
<i>X</i>	<i>Y</i>				
50	48,8	48,4	48,6	48,7	48,2
80	46,3	46,3	46,2	46,2	
90	45,6	45,3	45,8		
110	44,9	44,8	44,9	44,7	
130	43,3	43,4	43,4		
160	41,8	41,5	42,0		
170	41,5	41,5	41,7	41,4	
180	40,5	40,6	40,6		
200	39,9	39,8	40,2	40,4	
230	38,9	38,8	38,7	39,0	38,6

Продолжение таблицы №2

Вариант 4					
<i>X</i>	<i>Y</i>				
0,10	5,33	5,32	5,31		
0,20	5,70	5,69	5,73	5,73	5,71
0,25	5,93	5,91	5,87	5,94	
0,30	6,07	6,08	6,12	6,09	6,09
0,40	6,38	6,38	6,41		
0,65	7,13	7,15	7,09		
0,80	7,43	7,44	7,44	7,45	7,44
0,85	7,56	7,52	7,53	7,58	
0,90	7,63	7,66	7,62	7,59	7,61
1,05	7,82	7,86	7,86		
Вариант 5					
<i>X</i>	<i>Y</i>				
0	37	37	37	39	37
10	60	63	63		
20	81	83	82		
30	94	97	97	97	96
50	115	115	112		
60	121	118	118		
70	117	118	118	121	118
75	115	117	114		
85	110	111	111		
100	90	90	89	93	88
Вариант 6					
<i>X</i>	<i>Y</i>				
10	50,6	50,7	50,7		
30	49,2	48,8	49,4	49,1	
50	47,2	47,8	47,9	48,1	47,7
70	46,8	46,2	46,4	46,1	46,6
100	44,4	44,5	44,8		
110	43,9	43,9	43,9		
150	42,1	41,8	41,8	41,7	41,6
170	40,7	41,0	40,7	41,0	40,5
190	40,3	39,7	39,8	39,9	
220	38,4	38,7	38,8		

Продолжение таблицы №2

Вариант 7					
<i>X</i>	<i>Y</i>				
0,00	4,96	4,95	5,00	4,94	
0,05	5,21	5,23	5,23		
0,20	5,83	5,81	5,82	5,88	5,86
0,30	6,28	6,20	6,22		
0,40	6,59	6,60	6,55	6,54	6,54
0,50	6,86	6,89	6,90	6,88	6,87
0,65	7,32	7,34	7,26		
0,70	7,42	7,39	7,36	7,35	7,44
0,80	7,65	7,64	7,63		
0,90	7,81	7,77	7,85	7,90	
Вариант 8					
<i>X</i>	<i>Y</i>				
5	52	51	52	52	52
25	90	90	88		
35	102	104	102	105	103
40	110	109	106	106	
45	111	113	112		
50	114	115	112		
60	119	119	117	117	
65	117	118	118	120	116
80	112	115	110		
95	98	97	97	100	97
Вариант 9					
<i>X</i>	<i>Y</i>				
10	51,4	51,4	50,9		
40	49,7	49,5	49,2	49,1	49,0
50	48,2	48,7	48,7		
60	47,6	47,8	47,7		
100	45,0	45,3	45,3	45,1	44,9
120	43,6	43,7	43,8	43,9	44,1
150	42,2	42,1	42,3		
180	40,9	40,9	41,0	40,6	41,1
190	40,8	40,6	40,7		
200	40,1	40,1	40,2		

Окончание таблицы №2					
Вариант 10					
<i>X</i>	<i>Y</i>				
0,00	4,89	4,96	4,88	4,89	4,89
0,05	5,13	5,15	5,13	5,15	
0,20	5,70	5,76	5,71		
0,30	6,06	6,06	6,05	6,07	
0,40	6,40	6,40	6,38		
0,55	6,87	6,86	6,86		
0,60	7,00	7,01	6,96	7,02	
0,65	7,10	7,10	7,11		
0,85	7,55	7,49	7,56	7,54	
0,90	7,68	7,64	7,62	7,63	7,63

Некоторые теоретические сведения

1. Для построения линейной регрессионной модели $y = ax + b$ по методу наименьших квадратов коэффициенты a и b ищут исходя из условия минимизации взвешенной суммы квадратов отклонений результатов эксперимента от точек на прямой: $MinZ = \sum_{i=1}^n (ax_i + b - y_i)^2 w_i$. Нахождение минимума этой функции после нахождения частных производных по параметрам a и b сведется к решению системы линейных уравнений:

$$\begin{cases} a \sum_{i=1}^n x_i^2 w_i + b \sum_{i=1}^n x_i w_i = \sum_{i=1}^n x_i y_i w_i \\ a \sum_{i=1}^n x_i w_i + b \sum_{i=1}^n w_i = \sum_{i=1}^n y_i w_i \end{cases}$$

2. Для построения квадратичной регрессионной модели $y = ax^2 + bx + c$ по методу наименьших квадратов коэффициенты a , b и c ищут исходя из условия минимизации взвешенной суммы квадратов отклонений результатов эксперимента от точек на параболы: $MinZ = \sum_{i=1}^n (ax_i^2 + bx_i + c - y_i)^2 w_i$.

Нахождение минимума этой функции после нахождения частных производных по параметрам a , b и c сведется к решению системы линейных уравнений:

$$\begin{cases} a \sum_{i=1}^n x_i^4 w_i + b \sum_{i=1}^n x_i^3 w_i + c \sum_{i=1}^n x_i^2 w_i = \sum_{i=1}^n x_i^2 y_i w_i \\ a \sum_{i=1}^n x_i^3 w_i + b \sum_{i=1}^n x_i^2 w_i + c \sum_{i=1}^n x_i w_i = \sum_{i=1}^n x_i y_i w_i \\ a \sum_{i=1}^n x_i^2 w_i + b \sum_{i=1}^n x_i w_i + c \sum_{i=1}^n w_i = \sum_{i=1}^n y_i w_i \end{cases} .$$

3. Аналогично строятся и другие регрессионные модели.

4. Заметим, что из рассматриваемых программ только StatGraph позволяет сразу искать регрессионные модели с учетом весов.

5. Для проверки адекватности построенных моделей в том случае, если имел место повторный эксперимент, используют следующую схему.

Вычисляют $s_{ad}^2 = \sum_{i=1}^n \frac{(\Delta Y_i)^2 w_i}{k_{ad}}$, где $\Delta Y_i = (ax_i + b - y_i)$ для линейной модели (a и

b – найденные коэффициенты регрессионной модели), $\Delta Y_i = ax_i^2 + bx_i + c - y_i$ для квадратичной модели (a , b и c – найденные коэффициенты регрессионной модели), $k_{ad} = n - l$, где l – число оцениваемых параметров, то есть $l = 2$ для линейной модели, $l = 3$ – для квадратичной. Величина s_{ad}^2 (дисперсия адекватности) иногда обозначается $s_{ост}^2$ (остаточная дисперсия). Затем

вычисляют $s_{св}^2 = \frac{\sum_{i=1}^n s_i^2 k_i}{\sum_{i=1}^n k_i}$, где $s_i^2 = \sum_{j=1}^{w_i} \frac{(y_{ij} - \bar{y}_i)^2}{k_i}$, $k_i = w_i - 1$. Если $s_{ad}^2 < s_{св}^2$, то

модель является адекватной. Если нет, то находят $F_9 = \frac{s_{ad}^2}{s_{св}^2}$ и сравнивают с

квантилью распределения Фишера $F_m = F(p, k_{ad}, k_{св})$, где $k_{св} = \sum_{i=1}^n k_i$, $p = 1 - \alpha$,

а α – уровень значимости. Если $F_9 < F_m$, то модель является адекватной на уровне значимости α , в противном случае – нет.

6. Для проверки адекватности построенных моделей в том случае, если не имел место повторный эксперимент, используют следующую схему. Сначала вычисляют s_{ad}^2 (см. пункт 5), все веса принимаются равными 1. Затем

вычисляют $s_{mod}^2 = \sum_{i=1}^n \frac{(\Delta \bar{Y}_i)^2}{l-1}$, где $\Delta \bar{Y}_i = ax_i + b - \bar{y}$ для линейной модели,

$\Delta \bar{Y}_i = ax_i^2 + bx_i + c - \bar{y}$ – для квадратичной, а $\bar{y} = \sum_{i=1}^n \frac{y_i}{n}$. Величина s_{mod}^2

(дисперсия модели) иногда обозначается $s_{рег}^2$ (регрессионная дисперсия). После

этого сравнивают величину $F_9 = \frac{s_{mod}^2}{s_{ad}^2}$ с квантилью распределения Фишера

$F_m = F(p, l-1, n-l)$. Если окажется, что $F_9 < F_m$, то гипотеза об адекватности модели экспериментальным данным отвергается (на уровне значимости α), в противном случае гипотеза принимается.

7. Коэффициент уравнения регрессии называется незначимым, если его математическое ожидание (истинное значение) равно нулю. Для проверки значимости найденных коэффициентов регрессионной модели находят отношения этих коэффициентов к их стандартным ошибкам. Если через A обозначить матрицу системы для нахождения коэффициентов регрессионной модели, то стандартную ошибку i -го коэффициента d_i можно найти так:

$$s(d_i) = \sqrt{s_{ad}^2 \cdot b_{ii}},$$

где b_{ii} – соответствующий элемент матрицы $B = A^{-1}$.

Если модуль этого отношения больше квантили модуля отношения Стьюдента для числа степеней свободы k_{ad} , то есть $\left| \frac{d_i}{s(d_i)} \right| > |t|_p$ ($k = k_{ad}$), то коэффициент считается значимым, в противном случае – нет.

Пояснения к работе с программой MathCad (задание 1)

1. Следует подготовить два текстовых файла, например, в «Блокноте» со значениями x и средних арифметических y .

2. Укажем некоторые команды, которые нам понадобятся для выполнения работы.

slope(x,y)

Возвращает коэффициент при x в линейной регрессионной модели.

intercept(x,y)

Возвращает свободный член в линейной регрессионной модели.

regress(x,y,l-1)

Возвращает вектор, последние l координат которого – коэффициенты степенной регрессионной модели степени $l-1$.

3. MathCad позволяет искать практически любые регрессионные модели, даже не линейные относительно параметров. Для этого необходимо задать вектор, координатами которого являются приближающая функция и ее частные производные по параметрам (в примере – это вектор F). Затем следует задать начальное приближение для параметров (в примере – это вектор vb). Начальное приближение можно брать наугад, но затем всегда следует проверять правильность построенной модели графически, так как в случае неудачно подобранного начального приближения программа может выдать неправильный результат. Далее надо воспользоваться командой **genfit(x,y,vb,F)**, которая возвратит коэффициенты регрессионной модели.

4. Индексация в массивах в MathCad'е начинается с нуля.

Пояснения к работе с программой Excel

1. Для построения линейной регрессионной модели в программе Excel следует воспользоваться командой **ЛИНЕЙН(Y,X)**. И так как результатом этой команды является массив (из двух чисел), то, указав диапазон ячеек для x и для y , следует нажать комбинацию клавиш Ctrl+Shift+Enter.

2. Если же указать диапазон ячеек из двух столбцов и пяти строк, то программа выдаст дополнительную статистику, а именно:

Коэффициент a линейной модели	Коэффициент b линейной модели
Стандартная ошибка коэф-та a	Стандартная ошибка коэф-та b
Коэффициент детерминации	Стандартная ошибка для y (s_{ad})
F -наблюдаемое (F_3)	Степени свободы ($k_{ad} = n - l$)
Регрессионная сумма квадратов $((\Delta \bar{Y}_i)^2)$	Остаточная сумма квадратов $((\Delta Y_i)^2)$

Пояснения к работе с программой StatGraph

1. Для построения линейной регрессионной модели воспользуйтесь пунктом меню «*Relate / Multiple Regression...*», для построения квадратичной модели – «*Relate / Polynomial Regression...*».

2. Столбец «*p-Value*» из первой таблицы показывает уровни значимости коэффициентов регрессионной модели. Если уровень значимости некоторого коэффициента не превосходит заданного уровня значимости α , то этот коэффициент считается значимым, в противном случае – незначимым. Столбец «*p-Value*» из второй таблицы показывает уровень значимости F_3 . Если уровень

значимости F_{α} не превосходит заданного уровня значимости α , то модель считается адекватной, в противном случае – не адекватной.

Пояснения к работе с программой Stadia

1. Следует подготовить два столбца с экспериментальными данными, один – со значениями x , второй – со средними арифметическими для y .

2. В качестве десятичного разделителя следует использовать точку.

3. Для выполнения задания воспользуемся пунктом меню «Статист- $F9$ ». В появившемся окне следует нажать на кнопку «Простая регрессия/тренд», затем выбрать вид регрессионной модели.

4. Строка «Значим.» из первой таблицы показывает уровни значимости коэффициентов регрессионной модели. Вывод относительно адекватности построенной модели выдается текстовым сообщением.

Пояснения к работе с программой MathCad (задание 5)

1. Для выполнения работы понадобятся панели инструментов «Арифметика», «Матанализ», «Программирование», «Матрицы», «Булево». Вывести их на экран можно через пункт меню «Вид / Панель инструментов».

2. Вновь будем составлять по возможности универсальную программу для своего задания.

3. Перед тем, как приступить к работе с MathCad'ом, следует подготовить два текстовых файла, например, с помощью программы Блокнот. Один – со значениями x (в нашем примере эти значения записаны в столбец), второй – со значениями y (в нашем примере для фиксированного x соответствующие значения y записаны в строки). В качестве десятичного разделителя следует использовать точку. Сохраните эти текстовые файлы в той же папке, где будет храниться файл, созданный в MathCad'е.

4. Для выполнения работы нам понадобятся кнопки «Add Line», «Локальное присвоение», «Цикл For», «Оператор If» (панель «Программирование»); «Суммирование» (панель «Матанализ»); «Булево равенство» (панель «Булево»); «Создать матрицу или вектор», «Нижний индекс» (панель «Матрицы»).

5. Команда **Isolve(a,b)** возвращает вектор, являющийся решением системы линейных уравнений с матрицей системы a и вектором свободных членов b .

6. Напомним, что индексация в массивах в MathCad'е начинается с нуля.

Разбор нулевого варианта

Таблица 3 – Содержание нулевого варианта к лабораторной работе № 3

<i>Нулевой вариант</i>					
<i>X</i>	<i>Y</i>				
10	50,7	50,6	50,6		
20	49,9	50,0	50,1	50,2	
40	48,3	48,4	48,4		
60	47,4	47,4	46,8	47,4	
80	46,0	45,9	46,2	46,1	45,5
120	43,2	43,3	43,1	43,5	43,0
140	41,9	42,2	41,9	42,2	
150	41,7	41,7	42,0		
180	40,2	40,4	40,5	40,1	
200	39,5	39,9	39,7		

Программа MathCad (задание 1)

Разбор нулевого варианта

Подготовка

$X := \text{READPRN}(\text{"Var0-x.txt"})$ $Y := \text{READPRN}(\text{"Var0-ys.txt"})$

Линейное уравнение регрессии

$\text{slope}(X, Y) = -0.060$ $\text{intercept}(X, Y) = 50.886$

Ответ $y = -0.060x + 50.886$

Квадратичное и кубическое уравнения

$Z1 := \text{regress}(X, Y, 2)$ $Z2 := \text{regress}(X, Y, 3)$

$$Z1 = \begin{pmatrix} 3 \\ 3 \\ 2 \\ 51.566 \\ -0.081 \\ 1.073 \times 10^{-4} \end{pmatrix} \quad Z2 = \begin{pmatrix} 3 \\ 3 \\ 3 \\ 51.316 \\ -0.067 \\ -6.525 \times 10^{-5} \\ 5.453 \times 10^{-7} \end{pmatrix}$$

Ответ $y = 0.0001073x^2 - 0.081x + 51.566$

Ответ $y = 0.0000005453x^3 - 0.00006525x^2 - 0.067x + 51.316$

X =

	0
0	10
1	20
2	40
3	60
4	80
5	120
6	140
7	150
8	180
9	200

Y =

	0
0	50.633
1	50.05
2	48.367
3	47.25
4	45.94
5	43.22
6	42.05
7	41.8
8	40.3
9	39.7

Экспоненциальное уравнение регрессии, график

$$F(w, u) := \begin{pmatrix} e^{u_0 \cdot w + u_1} \\ w \cdot e^{u_0 \cdot w + u_1} \\ e^{u_0 \cdot w + u_1} \end{pmatrix} \quad vb := \begin{pmatrix} 0.1 \\ 1 \end{pmatrix} \quad p := \text{genfit}(X, Y, vb, F) \quad p = \begin{pmatrix} -1.336 \times 10^{-3} \\ 3.935 \end{pmatrix}$$

$g(r) := F(r, p)_0$

$i := 0..9$

Ответ $y = e^{(-0.001336x + 3.935)}$

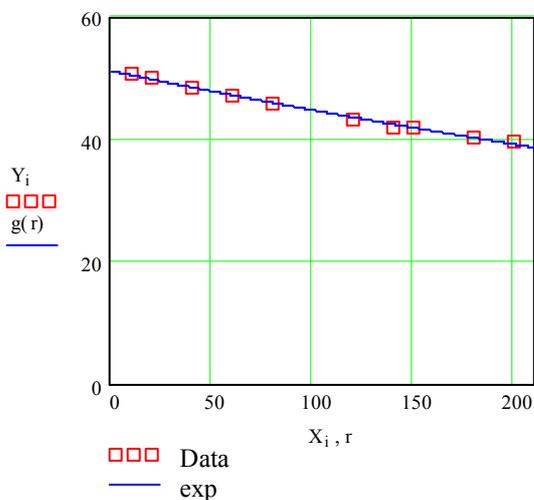
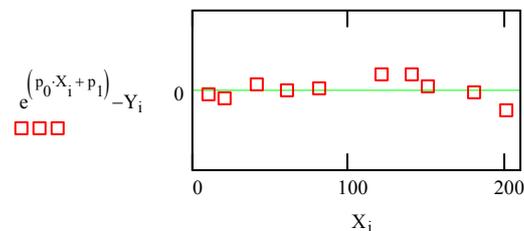


График остатков



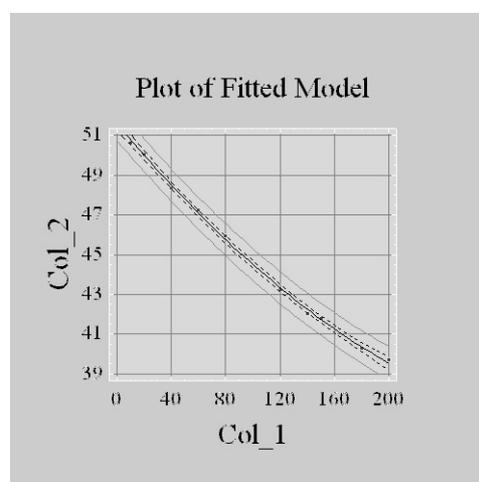
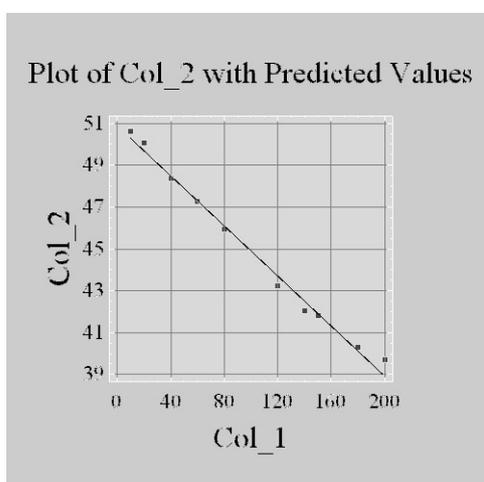
Программа Excel

Разбор нулевого варианта

X	Y				
10	50,633333				
20	50,05				
40	48,366667				
60	47,25				
80	45,94				
120	43,22				
140	42,05				
150	41,8				
180	40,3				
200	39,7				
Линейная модель:					
-0,0595463	50,8856341	Коэффициенты	-28,984815	208,593654	Отношения коэф-в к их ошибкам
0,0020544	0,24394622	Их стандартные ошибки	2,30600563	 t (k=8)	Оба коэф. знач., т.к. 29,0>2,3 и 208,6>2,3
0,99056737	0,41598385	Коэф. детерм. / Станд. ош. для y			
840,119521		8 F-наблюдаемое / Степени свободы	5,31764499	F(1;8)	Модель адекватна, т.к. 5,3<840,1
145,376437	1,38434052	Регрес. / Остаточ. суммы квадратов			

Программа StatGraph

Разбор нулевого варианта



Multiple Regression Analysis

 Dependent variable: Col_2

Parameter	Estimate	Standard Error	T Statistic	P-Value
CONSTANT	50,8814	0,248811	204,498	0,0000
Col_1	-0,0598933	0,00212574	-28,1753	0,0000

 Analysis of Variance

Source	Sum of Squares	Df	Mean Square	F-Ratio	P-Value
Model	504,362	1	504,362	793,85	0,0000
Residual	5,0827	8	0,635337		
Total (Corr.)	509,444	9			

R-squared = 99,0023 percent
 R-squared (adjusted for d.f.) = 98,8776 percent
 Standard Error of Est. = 0,79708
 Mean absolute error = 0,300386
 Durbin-Watson statistic = 0,658416

The StatAdvisor

The output shows the results of fitting a multiple linear regression model to describe the relationship between Col_2 and 1 independent variables. The equation of the fitted model is

$$\text{Col}_2 = 50,8814 - 0,0598933 \cdot \text{Col}_1$$

Since the P-value in the ANOVA table is less than 0.01, there is a statistically significant relationship between the variables at the 99% confidence level.

The R-Squared statistic indicates that the model as fitted explains 99,0023% of the variability in Col_2. The adjusted R-squared statistic, which is more suitable for comparing models with different numbers of independent variables, is 98,8776%. The standard error of the estimate shows the standard deviation of the residuals to be 0,79708. This value can be used to construct prediction limits for new observations by selecting the Reports option from the text menu. The mean absolute error (MAE) of 0,300386 is the average value of the residuals. The Durbin-Watson (DW) statistic tests the residuals to determine if there is any significant correlation based on the order in which they occur in your data file. Since the DW value is less than 1.4, there may be some indication of serial correlation. Plot the residuals versus row order to see if there is any pattern which can be seen.

In determining whether the model can be simplified, notice that the highest P-value on the independent variables is 0,0000, belonging to Col_1. Since the P-value is less than 0.01, the highest order term is statistically significant at the 99% confidence level. Consequently, you probably don't want to remove any variables from the model.

Polynomial Regression Analysis

 Dependent variable: Col_2

Parameter	Estimate	Standard Error	T Statistic	P-Value
CONSTANT	51,5913	0,154053	334,894	0,0000
Col_1	-0,0815422	0,00358172	-22,7662	0,0000
Col_1^2	0,000106205	0,0000170235	6,23871	0,0004

 Analysis of Variance

Source	Sum of Squares	Df	Mean Square	F-Ratio	P-Value
Model	508,67	2	254,335	2297,88	0,0000
Residual	0,774776	7	0,110682		
Total (Corr.)	509,444	9			

R-squared = 99,8479 percent
 R-squared (adjusted for d.f.) = 99,8045 percent
 Standard Error of Est. = 0,33269
 Mean absolute error = 0,25194
 Durbin-Watson statistic = 1,9561

The StatAdvisor

The output shows the results of fitting a second order polynomial model to describe the relationship between Col_2 and Col_1. The equation of the fitted model is

$$\text{Col}_2 = 51,5913 - 0,0815422 * \text{Col}_1 + 0,000106205 * \text{Col}_1^2$$

Since the P-value in the ANOVA table is less than 0.01, there is a statistically significant relationship between Col_2 and Col_1 at the 99% confidence level.

The R-Squared statistic indicates that the model as fitted explains 99,8479% of the variability in Col_2. The adjusted R-squared statistic, which is more suitable for comparing models with different numbers of independent variables, is 99,8045%. The standard error of the estimate shows the standard deviation of the residuals to be 0,33269. This value can be used to construct prediction limits for new observations by selecting the Forecasts option from the text menu. The mean absolute error (MAE) of 0,25194 is the average value of the residuals. The Durbin-Watson (DW) statistic tests the residuals to determine if there is any significant correlation based on the order in which they occur in your data file. Since the DW value is greater than 1.4, there is probably not any serious autocorrelation in the residuals.

In determining whether the order of the polynomial is appropriate, note first that the P-value on the highest order term of the polynomial equals 0,000428817. Since the P-value is less than 0.01, the highest order term is statistically significant at the 99% confidence level. Consequently, you probably don't want to consider any model of lower order.

Программа Stadia

Разбор нулевого варианта

ПРОСТАЯ РЕГРЕССИЯ. Файл var0.std

Переменные: x, y

Модель: линейная $Y = a_0 + a_1 \cdot x$

Коэфф.	a0	a1
Значение	50,89	-0,05955
Ст.ошиб.	0,2439	0,002054
Значим.	0	0

Источник	Сум.квадр.	Степ.св	Средн.квадр.
Регресс.	145,4	1	145,4
Остаточн	1,384	8	0,173
Вся	146,8	9	

Множеств R	R^2	R^2прив	Ст.ошиб.	F	Значим
0,99527	0,99057	0,98939	0,41598	840,1	0

Гипотеза 1: <Регрессионная модель адекватна экспериментальным данным>

Модель: парабола $Y = a_0 + a_1 \cdot x + a_2 \cdot x^2$

Коэфф.	a0	a1	a2
Значение	51,57	-0,08147	0,0001073
Ст.ошиб.	0,1431	0,003461	1,644E-5
Значим.	0	0	0,0006

Источник	Сум.квадр.	Степ.св	Средн.квадр.
Регресс.	146,6	2	73,28
Остаточн	0,1955	7	0,02793
Вся	146,8	9	

Множеств R	R^2	R^2прив	Ст.ошиб.	F	Значим
0,99933	0,99867	0,99829	0,16711	2624	0

Гипотеза 1: <Регрессионная модель адекватна экспериментальным данным>

Модель: полином $Y = \text{сумма}\{a_i \cdot x^i\}$

Коэфф.	a0	a1	a2	a3
Значение	51,32	-0,0669	-6,525E-5	5,453E-7
Ст.ошиб.	0,1635	0,00729	8,081E-5	2,519E-7
Значим.	0	0,0003	0,5453	0,072

Источник	Сум.квадр.	Степ.св	Средн.квадр.
Регресс.	146,7	3	48,88
Остаточн	0,1098	6	0,01829
Вся	146,8	9	

Множеств R	R^2	R^2прив	Ст.ошиб.	F	Значим
0,99963	0,99925	0,99888	0,13526	2672	0

Гипотеза 1: <Регрессионная модель адекватна экспериментальным данным>

Модель: экспонента $Y = \text{EXP}(a_0 + a_1 \cdot x)$

Коэфф.	a0	a1
Значение	3,934	-0,001326
Ст.ошиб.	0,004036	3,399E-5
Значим.	0	0

Источник	Сум.квадр.	Степ.св	Средн.квадр.
Регресс.	0,07213	1	0,07213
Остаточн	0,0003788	8	4,736E-5
Вся	0,07251	9	

Множеств R	R^2	R^2прив	Ст.ошиб.	F	Значим
0,99738	0,99478	0,99412	0,0068816	1523	0

Гипотеза 1: <Регрессионная модель адекватна экспериментальным данным>

Программа Mathcad (задание 5)

Разбор нулевого варианта

Подготовка

$x := \text{READPRN}(\text{"Var0-x.txt"})$

$y := \text{READPRN}(\text{"Var0-y.txt"})$

$n := 10 \quad p := 0.95$

$\text{ysr} := \left[\begin{array}{l} \text{for } i \in 0..n-1 \\ \text{ysr}_i \leftarrow \sum_{j=0}^{w_i-1} \frac{y_{i,j}}{w_i} \end{array} \right]$

$w := \begin{pmatrix} 3 \\ 4 \\ 3 \\ 4 \\ 5 \\ 5 \\ 4 \\ 3 \\ 4 \\ 3 \end{pmatrix}$

$x =$

	0
0	10
1	20
2	40
3	60
4	80
5	120
6	140
7	150
8	180
9	200

$y =$

	0	1	2	3	4
0	50.7	50.6	50.6	0	0
1	49.9	50	50.1	50.2	0
2	48.3	48.4	48.4	0	0
3	47.4	47.4	46.8	47.4	0
4	46	45.9	46.2	46.1	45.5
5	43.2	43.3	43.1	43.5	43
6	41.9	42.2	41.9	42.2	0
7	41.7	41.7	42	0	0
8	40.2	40.4	40.5	40.1	0
9	39.5	39.9	39.7	0	0

Построение линейной модели

$$a1 := \begin{bmatrix} \sum_{i=0}^{n-1} (x_i)^2 \cdot w_i & \sum_{i=0}^{n-1} x_i \cdot w_i \\ \sum_{i=0}^{n-1} x_i \cdot w_i & \sum_{i=0}^{n-1} w_i \end{bmatrix}$$

$$a2 := \begin{pmatrix} \sum_{i=0}^{n-1} x_i \cdot w_i \cdot \text{ysr}_i \\ \sum_{i=0}^{n-1} \text{ysr}_i \cdot w_i \end{pmatrix}$$

$a := \text{lsolve}(a1, a2)$

$$a = \begin{pmatrix} -0.060 \\ 50.881 \end{pmatrix}$$

$b := a1^{-1}$

Ответ: $y = -0.060x + 50.881$

Построение квадратичной модели

$$A1 := \begin{bmatrix} \sum_{i=0}^{n-1} (x_i)^4 \cdot w_i & \sum_{i=0}^{n-1} (x_i)^3 \cdot w_i & \sum_{i=0}^{n-1} (x_i)^2 \cdot w_i \\ \sum_{i=0}^{n-1} (x_i)^3 \cdot w_i & \sum_{i=0}^{n-1} (x_i)^2 \cdot w_i & \sum_{i=0}^{n-1} (x_i) \cdot w_i \\ \sum_{i=0}^{n-1} (x_i)^2 \cdot w_i & \sum_{i=0}^{n-1} (x_i) \cdot w_i & \sum_{i=0}^{n-1} w_i \end{bmatrix}$$

$$A2 := \begin{bmatrix} \sum_{i=0}^{n-1} (x_i)^2 \cdot w_i \cdot \text{ysr}_i \\ \sum_{i=0}^{n-1} x_i \cdot w_i \cdot \text{ysr}_i \\ \sum_{i=0}^{n-1} w_i \cdot \text{ysr}_i \end{bmatrix}$$

$A := \text{lsolve}(A1, A2)$

$$A = \begin{pmatrix} 1.062 \times 10^{-4} \\ -0.082 \\ 51.591 \end{pmatrix}$$

$B := A1^{-1}$

Ответ: $y = 0.0001062x^2 - 0.082x + 51.591$

Проверка адекватности обеих моделей

$$Sad1 := \sum_{i=0}^{n-1} (a_0 \cdot x_i + a_1 - ysr_i)^2 \cdot w_i$$

$$Sad1 := \frac{1}{n-2} \cdot Sad1 \quad Sad1 = 0.6353$$

$$Sad2 := \sum_{i=0}^{n-1} [A_0 \cdot (x_i)^2 + A_1 \cdot x_i + A_2 - ysr_i]^2 \cdot w_i$$

$$Sad2 := \frac{1}{n-3} \cdot Sad2 \quad Sad2 = 0.1107$$

$$Ssv := \frac{\sum_{i=0}^{n-1} \left[\sum_{j=0}^{w_i-1} \frac{(y_{i,j} - ysr_i)^2}{w_i - 1} \right] \cdot (w_i - 1)}{\sum_{i=0}^{n-1} (w_i - 1)}$$

$$Ssv = 0.0394$$

$$F1 := \frac{Sad1}{Ssv} \quad F1 = 16.1234$$

$$F2 := \frac{Sad2}{Ssv} \quad F2 = 2.8089$$

$$Fkv1 := qF \left[p, n - 2, \sum_{i=0}^{n-1} (w_i - 1) \right]$$

$$Fkv2 := qF \left[p, n - 3, \sum_{i=0}^{n-1} (w_i - 1) \right]$$

$$Fkv1 = 2.291$$

$$Fkv2 = 2.359$$

$$FlagLi1 := \begin{cases} \text{"Yes"} & \text{if } (Sad1 < Ssv) \\ \text{"No"} & \text{if } (Sad1 > Ssv) \end{cases}$$

$$FlagSq1 := \begin{cases} \text{"Yes"} & \text{if } (Sad2 < Ssv) \\ \text{"No"} & \text{if } (Sad2 > Ssv) \end{cases}$$

$$FlagLi2 := \begin{cases} \text{"Yes"} & \text{if } F1 < Fkv1 \\ \text{"No"} & \text{if } F1 > Fkv1 \end{cases}$$

$$FlagSq2 := \begin{cases} \text{"Yes"} & \text{if } F2 < Fkv2 \\ \text{"No"} & \text{if } F2 > Fkv2 \end{cases}$$

$$FlagLi := \begin{cases} \text{"Yes"} & \text{if } (FlagLi1 = \text{"Yes"}) \\ FlagLi2 & \text{if } FlagLi1 \neq \text{"Yes"} \end{cases}$$

$$FlagSq := \begin{cases} \text{"Yes"} & \text{if } (FlagSq1 = \text{"Yes"}) \\ FlagSq2 & \text{if } FlagSq1 \neq \text{"Yes"} \end{cases}$$

$$FlagLi = \text{"No"}$$

$$FlagSq = \text{"No"}$$

Ответ: обе модели не адекватны

Проверка значимости коэффициентов

Стандартные ошибки коэффициентов

$$s_{11} := \sqrt{Sad1} \cdot \sqrt{b_{0,0}} \quad s_{11} = 2.126 \times 10^{-3}$$

$$s_{12} := \sqrt{Sad1} \cdot \sqrt{b_{1,1}} \quad s_{12} = 0.249$$

$$s_{21} := \sqrt{Sad2} \cdot \sqrt{B_{0,0}} \quad s_{21} = 1.702 \times 10^{-5}$$

$$s_{22} := \sqrt{Sad2} \cdot \sqrt{B_{1,1}} \quad s_{22} = 3.582 \times 10^{-3}$$

$$s_{23} := \sqrt{Sad2} \cdot \sqrt{B_{2,2}} \quad s_{23} = 0.154$$

Отношения коэффициентов к их ошибкам

$$t_{11} := \frac{a_0}{s_{11}} \quad t_{12} := \frac{a_1}{s_{12}} \quad t_{21} := \frac{A_0}{s_{21}} \quad t_{22} := \frac{A_1}{s_{22}} \quad t_{23} := \frac{A_2}{s_{23}}$$

Проверка значимости

$$T_{kv1} := -qt\left(\frac{1-p}{2}, n-2\right) \quad T_{kv1} = 2.306 \quad T_{kv2} := -qt\left(\frac{1-p}{2}, n-3\right) \quad T_{kv2} = 2.365$$

$$\text{Flag}_{11} := \begin{cases} \text{"a1_Znachim"} & \text{if } (|t_{11}| > T_{kv1}) \\ \text{"a1_Ne_Znachim"} & \text{if } (|t_{11}| < T_{kv1}) \end{cases} \quad \text{Flag}_{12} := \begin{cases} \text{"a2_Znachim"} & \text{if } (|t_{12}| > T_{kv1}) \\ \text{"a2_Ne_Znachim"} & \text{if } (|t_{12}| < T_{kv1}) \end{cases}$$

$$\text{Flag}_{21} := \begin{cases} \text{"A1_Znachim"} & \text{if } (|t_{21}| > T_{kv2}) \\ \text{"A1_Ne_Znachim"} & \text{if } (|t_{21}| < T_{kv2}) \end{cases} \quad \text{Flag}_{22} := \begin{cases} \text{"A2_Znachim"} & \text{if } (|t_{22}| > T_{kv2}) \\ \text{"A2_Ne_Znachim"} & \text{if } (|t_{22}| < T_{kv2}) \end{cases}$$

$$\text{Flag}_{23} := \begin{cases} \text{"A3_Znachim"} & \text{if } (|t_{23}| > T_{kv2}) \\ \text{"A3_Ne_Znachim"} & \text{if } (|t_{23}| < T_{kv2}) \end{cases}$$

$$\text{Flag}_{11} = \text{"a1_Znachim"}$$

$$\text{Flag}_{12} = \text{"a2_Znachim"}$$

$$\text{Flag}_{21} = \text{"A1_Znachim"}$$

$$\text{Flag}_{22} = \text{"A2_Znachim"}$$

$$\text{Flag}_{23} = \text{"A3_Znachim"}$$

Ответ: все коэффициенты значимы

Лабораторная работа №4: Исследование линейной корреляции

Используемое ПО: MathCad.

Цель работы: Научиться с помощью программы MathCad делать выводы о силе и характере связи между двумя величинами.

Задание

1. Найти эмпирический коэффициент корреляции. Найти уравнения эмпирических прямых регрессии y на x и x на y . На одном чертеже построить прямые регрессии.

2. Проверить гипотезу о незначимости коэффициента корреляции.

3. Построить доверительный интервал для коэффициента корреляции и сделать вывод о силе и характере связи между величинами.

Таблица 4 – Содержание вариантов к лабораторной работе №4

Вариант 1									
X	-3	1	2	-5	0	3	-1	5	7
Y	0,3	0,7	0,9	0,1	0,5	1,0	0,4	1,3	1,8
Вариант 2									
X	-3	1	2	-5	0	3	-1	5	7
Y	1,6	1,2	1,3	2,0	1,7	1,0	0,4	0,5	0,2
Вариант 3									
X	1,5	1,0	2,5	4,0	7,5	5,5	6,0	7,0	9,0
Y	0,27	0,15	0,34	0,42	0,91	0,58	0,72	0,84	1,12
Вариант 4									
X	1,5	1,0	2,5	4,0	7,5	5,5	6,0	7,0	9,0
Y	1,98	2,03	1,54	1,33	0,51	1,20	1,09	0,33	0,31
Вариант 5									
X	2,1	1,4	2,7	3,9	1,9	5,2	4,3	3,2	6,0
Y	3,6	0,1	3,5	5,1	2,0	9,1	7,4	4,3	9,8
Вариант 6									
X	2,1	1,4	2,7	3,9	1,9	5,2	4,3	3,2	6,0
Y	8,3	9,9	7,2	6,2	8,0	4,0	5,4	6,9	2,8
Вариант 7									
X	-3	1	2	-5	-4	9	-1	5	7
Y	0,4	1,7	1,9	0,3	0,5	3,6	0,9	2,3	2,8

Окончание таблицы №4									
Вариант 8									
<i>X</i>	-3	1	2	-5	-4	9	-1	5	7
<i>Y</i>	5,3	3,7	3,9	6,1	5,5	2,0	4,4	3,3	2,8
Вариант 9									
<i>X</i>	-3	1	-2	-5	6	3	-1	5	7
<i>Y</i>	0,23	0,47	0,39	0,12	2,35	1,04	0,42	2,23	3,02
Вариант 10									
<i>X</i>	-3	1	-2	-5	6	3	-1	5	7
<i>Y</i>	7,63	5,78	6,92	8,11	0,5	4,60	6,40	3,93	3,08

Некоторые теоретические сведения

1. Напомним, что эмпирические дисперсии рассчитываются по следующим формулам:

$$s_x^2 = \tilde{K}_{11} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}, \quad s_y^2 = \tilde{K}_{22} = \sum_{i=1}^n \frac{(y_i - \bar{y})^2}{n-1}.$$

Оценкой ковариации $\text{cov}(x, y) = K_{12}$ является величина

$$\tilde{K}_{12} = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1}.$$

Эмпирический коэффициент корреляции находится по формуле:

$$R = \frac{\tilde{K}_{12}}{\sqrt{\tilde{K}_{11} \tilde{K}_{22}}} = \frac{\tilde{K}_{12}}{s_x s_y}.$$

Уравнение эмпирической прямой регрессии *Y* на *X* задается так:

$$\frac{y - \bar{y}}{s_y} = R \frac{x - \bar{x}}{s_x};$$

уравнение прямой *X* на *Y* так:

$$\frac{y - \bar{y}}{s_y} = \frac{1}{R} \frac{x - \bar{x}}{s_x}.$$

2. Эмпирический коэффициент корреляции называется незначимым, если его истинное значение $\rho = 0$, то есть линейная зависимость между величинами *X* и *Y* отсутствует. Для проверки гипотезы о незначимости коэффициента корреляции сравним экспериментальное значение

$$T_{\vartheta} = \left| R \sqrt{\frac{n-2}{1-R^2}} \right|$$

с квантилью модуля отношения Стьюдента $|t|_p$ ($k = n - 2$). Если $T_{\vartheta} > |t|_p$, то гипотеза отклоняется, в противном случае гипотеза принимается.

3. Доверительный интервал для коэффициента корреляции имеет вид:

$$th\left(Z - \frac{|u|_p}{\sqrt{n-3}}\right) < \rho < th\left(Z + \frac{|u|_p}{\sqrt{n-3}}\right),$$

где $Z = \frac{1}{2} \ln\left(\frac{1+R}{1-R}\right)$, $|u|_p$ – квантиль модуля стандартного нормального распределения. Если $\rho = 0$ не принадлежит найденному доверительному интервалу, то гипотеза о существовании линейной зависимости принимается, как не противоречащая экспериментальным данным; в противном случае гипотеза о существовании линейной зависимости отвергается.

Пояснения к работе с программой MathCad

1. Для выполнения работы понадобятся панели инструментов «Арифметика», «Программирование», «Матрицы», «Булево», «Графики». Вывести их на экран можно через пункт меню «Вид / Панель инструментов».

2. Программу желательно составлять так, чтобы она обладала универсальностью. В том случае, если придется изменить выборки, то пусть в программе при этом практически ничего не придется менять.

3. Укажем некоторые новые команды, которые нам понадобятся для выполнения работы.

corr(x,y)

Эмпирический коэффициент корреляции.

qnorm(p,α,σ)

Обычная квантиль (не модуля) нормального распределения с параметрами (α, σ) для вероятности p . Заметим, что из обычной квантили можно получить нужную нам квантиль модуля так:

$$qnorm\left(\frac{1+p}{2}, \alpha, \sigma\right).$$

tanh(x)

Гиперболический тангенс числа x .

Разбор нулевого варианта

Таблица 5 – Содержание нулевого варианта к лабораторной работе № 4

<i>Нулевой вариант</i>									
<i>X</i>	-2,0	0,5	-1,5	-5,5	3,5	-1,0	2,0	0,0	1,5
<i>Y</i>	0,11	0,09	0,13	0,11	0,06	0,12	0,08	0,11	0,07

Программа MathCad Разбор нулевого варианта

0 Подготовка

```
X := ( -2 0.5 -1.5 -5.5 3.5 -1 2 0 1.5 )
Y := ( 0.11 0.09 0.13 0.11 0.06 0.12 0.08 0.11 0.07 )
n := 9
```

1 Эмпирический коэффициент корреляции, прямые регрессии

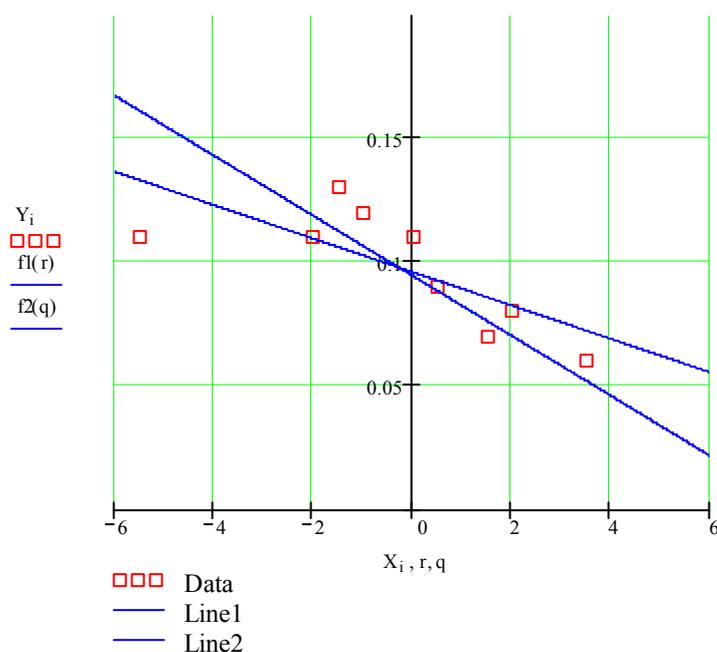
```
x1 := mean(X)           y1 := mean(Y)

sx := stdev(X) * sqrt(n / (n - 1))
sy := stdev(Y) * sqrt(n / (n - 1))

R := corr(X, Y)

f1(x) := y1 + R * (sy / sx) * (x - x1)
f2(x) := y1 + (1 / R) * (sy / sx) * (x - x1)

X := X^T               Y := Y^T               i := 0..8
```



2 Проверка гипотезы о незначимости коэффициента корреляции

$$t1 := \left| R \cdot \sqrt{\frac{n-2}{1-R^2}} \right| \quad t1 = 2.968 \quad t2 := -qt(0.025, n-2) \quad t2 = 2.365$$

$$\text{Flag} := \begin{cases} \text{"Znachim"} & \text{if } (t1 > t2) \\ \text{"Ne_Znachim"} & \text{if } (t2 \geq t1) \end{cases} \quad \text{Flag} = \text{"Znachim"}$$

3 Доверительный интервал для R

$$Z := 0.5 \cdot \ln\left(\frac{1+R}{1-R}\right) \quad Z = -0.965 \quad Sz := \sqrt{\frac{1}{n-3}} \quad Sz = 0.408$$

$$u := qnorm(0.975, 0, 1) \quad u = 1.96$$

$$r1 := \tanh(Z - u \cdot Sz) \quad r1 = -0.943 \quad r2 := \tanh(Z + u \cdot Sz) \quad r2 = -0.163$$

$$\text{Flag} := \begin{cases} \text{"Linei"} & \text{if } (r1 > 0) \vee (r2 < 0) \\ \text{"Ne_Linei"} & \text{if } (r1 < 0) \wedge (r2 > 0) \end{cases} \quad \text{Flag} = \text{"Linei"}$$

Библиографический список

1. Большев Л. Н. Таблицы математической статистики [Текст] : изд-е 3-е / Л. Н. Большев. – М. : Наука, 1983.
2. Вентцель Е.С. Теория вероятностей [Текст] : учебник / Е.С. Вентцель.– М. : Наука, 1969. – 576 с.
3. Гмурман В.Е. Теория вероятностей и математическая статистика [Текст] : изд-е 7-е, стер. / В.Е. Гмурман. – М. : Высш. шк., 2001. – 479с.
4. Калинина В.Н. Математическая статистика [Текст] : изд-е 3-е, испр. / В.Н.Калинина, В.Ф. Панкин. – М. :.Высш. шк., 2001. – 336 с.
5. Карасев В.А. Организация эксперимента [Текст] : пособие / В.А. Карасев, Л.З. Румшинский. – М. : ротاپринт МИСиС, 1986. – 86 с.
6. Кремер Н.Ш. Теория вероятностей и математическая статистика [Текст] : учебник , изд-е 3-е, перераб. и доп. / Н.Ш. Кремер. – М.: ЮНИТИ–ДАНА, 2007.– 551 с.
7. Румшинский Л.З. Организация эксперимента [Текст] : пособие / Л.З. Румшинский. – М. : ротاپринт МИСиС, 1984. – 140 с.
8. Тюрин Ю.Н. Анализ данных на компьютере [Текст] : учебник / Ю.Н. Тюрин, А.А. Макаров. – М. : Инфра-М, 2003. – 544 с.

ISBN

Учебное издание

Дмитрий Давидович Изаак
Анна Викторовна Швалёва

Математическая статистика

ЛАБОРАТОРНЫЙ ПРАКТИКУМ

*Работа отпечатана с оригинала-макета,
предоставленного авторами*

НФ НИТУ «МИСиС»

462359, Оренбургская обл., г. Новотроицк, ул. Фрунзе, дом 8
